

Assessing Fair Policing in Austin, TX

Team FunkyStats

4/18/2021

Abstract

This report demonstrates disparities by race in traffic stops by the Austin Police Department. After exploratory analysis, we assess various models and statistics derived from the hit rate and using the Veil of Darkness. We conclude with a Bayesian hierarchical model that produces officer-level posteriors for the hit rate.

Contents

Introduction	1
Available Data	1
Summary Statistics	1
Stanford Data	2
Exploratory Analysis	3
Benchmark Test	3
Outcome Test	4
Veil of Darkness Test & Fisher's Exact Test	4
Hit Rate and Causal Issues	5
Modeling	5
Logistic Regression	5
Logistic Regression for Frisk Rate	5
Logistic Regression for Contraband found	6
Bayesian Modeling	7
Investigating the Hit Rate	7
Conclusion	9
References	9

Introduction

This paper investigates racial disparities in traffic stops by the Austin Police Department. Using data available from Austin Open Data, a Texas government-run data portal, and from the Stanford Open Policing Project, we evaluate these disparities using models derived from the “hit rate” and the effect of the “veil of darkness,” two often-cited methods for assessing fair policing. Our main report consists of three parts. First, we conduct an exploratory data analysis to get a big picture of policing in Austin. Second, we use various modelling strategies to assess the severity of racial disparities. Third, we propose a measure of fairness based on the differences in the posterior median hit rate among individual police officers.

Available Data

The primary data set is from the Stanford Open Policing Project¹ (from here on referred to as the Stanford data). This data set record stops made by the APD a roughly ten year period (2006.01.01 - 2016.06.30) and contains information such as the date of the stops, the subject's race, whether the person was searched or frisked, whether any contraband were found. Notably, this data lacks information about the time or place of the stops. Because the 2016 data is incomplete, we focus on the data for which we have complete years (2006-2015) for the first two parts of the analysis; this contains 463,944 stops.

Our secondary data set is from the 2019 Racial Profiling report, available from Austin Open Data (and hereafter referred to as the RP data). This data set contains similar information as the Stanford data, with additional information about the event time, location, and officer race. Notably, the race of the subject is missing from this data.

Lastly, we use US census demographic data. Specifically, we use 2017 5-year American Community Survey zip-code-level data with the Stanford data, and 2019 census population data for 2019 Austin RP data. In addition, we also refer to the racial profiling reports from the Austin Police Department.

Summary Statistics

Stanford Data

Summary statistics for the Stanford data are as follows. The statistics reported cover all available data (2006.01.01 - 2016.06.30). Unique officer IDs are available but not shown here.

% latex table generated in R 4.0.5 by xtable 1.8-4 package % Sun Apr 18 22:42:00 2021

	nobs	nmis	uniq	mean	SD	min	25%	50%	75%	max
subject_age	480091	3164	94	37.98	13.82	10.00	26.00	36.00	48.00	103.00
subject_sex	482881	374	2	0.30	0.46	0.00	0.00	0.00	1.00	1.00
frisk_performed	483255	0	2	0.02	0.15	0.00	0.00	0.00	0.00	1.00
search_conducted	483255	0	2	0.04	0.20	0.00	0.00	0.00	0.00	1.00
search_person	483255	0	2	0.03	0.18	0.00	0.00	0.00	0.00	1.00
search_vehicle	483255	0	2	0.02	0.15	0.00	0.00	0.00	0.00	1.00

% latex table generated in R 4.0.5 by xtable 1.8-4 package % Sun Apr 18 22:42:10 2021

	nobs	nmis	uniq	mean	SD	min	25%	50%	75%	max
contraband_found	19256	0	2	0.25	0.43	0.00	0.00	0.00	0.00	1.00
contraband_drugs	19256	0	2	0.01	0.12	0.00	0.00	0.00	0.00	1.00
contraband_weapons	19256	0	2	0.05	0.21	0.00	0.00	0.00	0.00	1.00
frisk_performed	19256	0	2	0.51	0.50	0.00	0.00	1.00	1.00	1.00

¹Stanford Open Policing Project (OPP): <https://openpolicing.stanford.edu/data/>

Table 1: Subject race.

Race	n	percent
asian/pacific islander	13167	0.0272466
black	72324	0.1496607
hispanic	123943	0.2564764
other	2626	0.0054340
unknown	3135	0.0064873
white	268058	0.5546950
NA	2	NA

Table 2: Search basis.

Search basis	n	percent
consent	3195	0.1659223
other	276	0.0143332
plain view	152	0.0078936
probable cause	15633	0.8118509
NA	463999	NA

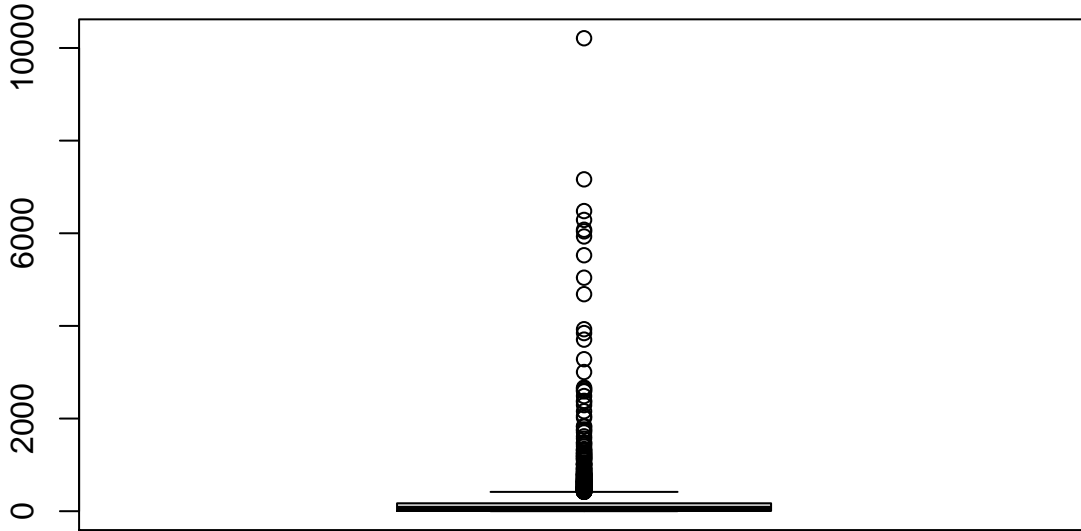


Figure 1: Distribution of stops by unique officer ID

We note that the distribution of stops per officer has an extremely long tail.

Exploratory Analysis

We first examine the count of stops by race during 2006-2015 (Table 3), using the Stanford Data. It is notable that over half of the stops involved were of white subjects, about four times the number of stops of Black people. According to 5-year census data, the white population in Austin (445,269) is almost 7 times than the black population (66,724) — a classic Simpson’s paradox. Examining figure 2, we can see that at least for Black, Hispanic and white drivers, the annual trends are very different by race.

Driver Race	Counts	Proportion
asian/pacific	11658	0.033
black	52381	0.147
hispanic	765707	0.215
other	2105	0.006
unknown	2622	0.007
white	211588	0.593

Table 3: Proportion of stops by race during 2006-2015

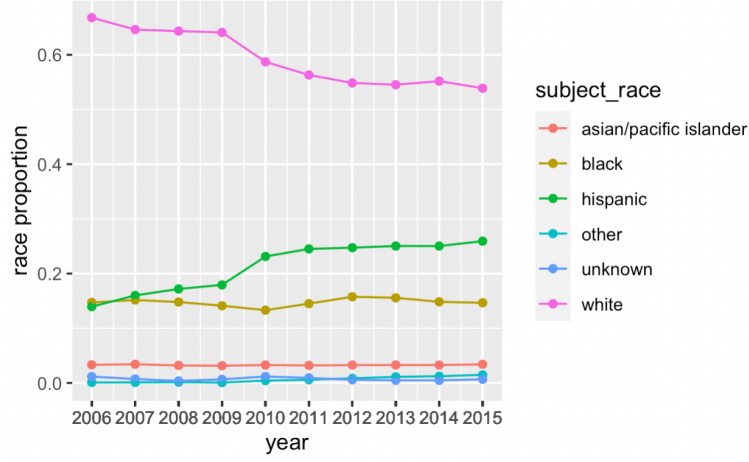


Figure 2: Race proportion in each year

We can see that fewer white drivers stopped especially after 2009, whereas there continued to be an increase trend for Hispanic and black drivers stopped over the nine years.

Benchmark Test

As mentioned previously, the number of the white drivers stopped is four times of the black drivers, while the white population in Austin is actually over 6 times of the black population. We need to compare the stop times with the population in each demographic groups.

$$\text{Stop Rate}_i = \frac{\text{Number of Stops for Race } i}{\text{Population of Race } i}$$

From Table 4, we can see black drivers are stopped with a rate much higher than the drivers of other races. We can investigate this further by looking at other benchmarks such as search rate and frisk rate with stopped population as baseline.

$$\text{Search Rate}_i = \frac{\text{Number of Stopped People Who Were Searched for Race } i}{\text{Number of Stops for Race } i}$$

$$\text{Frisk Rate}_i = \frac{\text{Number of Stopped People Who Were Frisked for Race } i}{\text{Number of Stops for Race } i}$$

Again, from the last two columns of Table 4, we can see Black and Hispanic drivers are searched with a rate 3 times higher than white drivers, and almost 8 times higher than Asian/Pacific drivers. The Black and Hispanic drivers are also frisked with a rate much higher than the drivers of other races.

Table 4: Stop rates, search rates and frisk rates during 2015

Driver Race	Counts	Population	Proportion	Stop Rate	Search Rate	Frisk Rate	Hit Rate
asian/pacific	11658	63752	0.033	0.183	0.015	0.011	0.188
black	52381	66724	0.147	0.785	0.092	0.039	0.254
hispanic	765707	316709	0.215	0.242	0.086	0.044	0.323
white	211588	445269	0.593	0.475	0.031	0.021	0.318

The racial disparity by the police is clear from benchmark test, but it is insufficient evidence of discriminative policing. The key part of this analysis is to find out the true distribution of the drivers violating the traffic laws or conducting crimes. We need to check if different race groups are disproportionately stopped corresponding to their rates of violating the law.

Outcome Test

In order to more rigorously investigate and measure the fairness, we shall look at the result of the searches to see if the targeted drivers are really actually doing something illegal. Here, we define a successful search as one that uncovers contraband, and we define the hit rate as the proportion of searches that are successful.

$$\text{Hit Rate}_i = \frac{\text{Number of Contraband Uncovered for Race } i}{\text{Number of Searched People for Race } i}$$

If racial groups have different hit rates, it can be taken as evidence of discriminative policing. From the last column in Table 4, we can see the hit rate for Black and Asian/Pacific drivers are lower than for white and Hispanic drivers, indicating police may have a lower threshold of evidence when searching Black or Asian/Pacific drivers.

Veil of Darkness Test & Fisher’s Exact Test

According to Grogger and Ridgeway, the “Veil of Darkness” test can help assess the bias in the stop decisions. The hypothesis of this test states that officers who are engaged in racial profiling are less likely to be able to identify a driver’s race after dark than during daylight. Under this hypothesis, if stops made after dark had smaller proportion of black drivers stopped than stops made during daylight, it could be evidence of racial profiling.

Because neither of our data sets contain both driver races and stop time information, we venture to an indirect way by measuring the racial population in different areas through zip codes by 2019 RP data provided by Austin Police Department. In order to accurately distinguish the daytime and nighttime, we compute the daily subset and dusk time for Austin in 2019. In Table 5 we can see earliest sunset in 2019 was at around 17:32 in early December and it goes fully dark in 26 minutes. The latest sunset time was around 20:38 late June and it was fully dark after 28 minutes.

Date	Sunset	Dusk	Sunset Minute	Dusk Minute
2019-12-02	17:31:42	17:57:48	1051	1077
2019-12-01	17:31:45	17:57:48	1051	1077
2019-06-30	20:37:58	21:05:27	1237	1265
2019-06-29	20:37:56	21:05:27	1237	1265

Table 5: Minimum and maximum dusk time during the 2019 in Austin

	Day	Night
BDA	124	126
WDA	2937	2216

Table 6: Contingency Table

We denote the stops happening before sunset as Daytime Stop, and the stop happening after the dusk as Nighttime Stop. We do not consider the stops happening between the sunset and dusk in this study. According to ZIP codes and the corresponding demographic data, we consider the areas that have more black people as black dominant area (BDA), and the areas consist of more white people as white dominated area (WDA). For simplicity of the analysis, here we consider only the black and the white population groups. Hence, each zip code is regarded as a location with label as white (WDA) or black (BDA).

We record the stops happening in each zip codes into two categories: daytime stops or nighttime stops, and we treat two rows in Table 6 as independent binomial samples. Of $n_1 = 250$ recorded stops in black dominated area, 124 stops happened during the daytime, a proportion of $p_1 = 124/250 = 0.496$. Of $n_2 = 5153$ recorded stops in white dominated area, 2937 stops happened during the daytime, a proportion of $p_2 = 2937/5153 = 0.570$. The sample difference of proportions is 0.074. We obtain Fisher’s exact test for testing null hypothesis of independence of the two rows with p value of 0.02, indicating the strong evidence that the police are not equally likely practicing during day and night to different racial groups.

Hit Rate and Causal Issues

In our analysis of observational data, we have to deal with unmeasured confounders. For example, one can argue that it would be unsurprising if more officers patrolled areas with higher crime rates; crime rates are known to be correlated with income and demographic factors. Therefore, if those neighborhoods have higher minority populations, we would expect to see more minority traffic stops. Although this still exposes problems in Austin, it could be interpreted as a problem of economic segregation, not traffic fairness.

To overcome this problem, we propose to look at the hit rate with more details in Section 3. We argue that given a person is being searched, the probability of finding contraband items should be equal among all races, regardless of the neighborhood that the search conducted. We want to emphasize that using hit rate does not eliminate all unmeasured confounders, but it helps mitigate the problem.

Modeling

Logistic Regression

Logistic Regression for Frisk Rate

Our descriptive analysis shows that black people in Austin seem to be more likely to be stopped by the police. We want to answer the question, given a person is stopped, what factors may impact the likelihood of that person being frisked? To investigate this, we fit a logistic regression model with **frisk** as the dependent variable and **race**, **age**, and **sex**.

$$\text{Logit}[P(\text{Being Frisked})] = \beta_0 + \beta_1 \text{Race} + \beta_2 \text{Age} + \beta_3 \text{Sex}$$

Results can be found below.

Logistic Regression for Contraband found

We want to investigate how likely contraband items are found when searching is performed. This is equivalent to calculating hit rate defined in section 2.2.2. We argue that if racial bias does not exist, the hit rate should be equal for all races. In other words, we expect to find that **race** is not an essential factor in the model:

$$\text{Logit}[P(\text{Contraband found})] = \beta_0 + \beta_1 \text{Race}.$$

We also break down contraband found into three categories: Drugs, Weapons, and Others. We also fit a logistic regression model for each of these categories with **Race** as the sole independent variable.

Table 7: Logistic model for frisk rate vs. race, age, and sex

term	estimate	std.error	statistic	p.value
(Intercept)	-2.984	0.102	-29.402	0.000
subject_raceblack	1.503	0.100	15.044	0.000
subject_racehispanic	1.310	0.099	13.214	0.000
subject_racewhite	0.719	0.099	7.260	0.000
subject_raceother	0.617	0.180	3.434	0.001
subject_raceunknown	0.647	0.183	3.544	0.000
subject_age	-0.046	0.001	-51.125	0.000
subject_sexfemale	-1.643	0.036	-45.311	0.000

Table 8: Logistic model for contraband found vs. race

term	estimate	std.error	statistic	p.value
(Intercept)	-1.988	0.222	-8.944	0.000
subject_raceblack	0.899	0.225	3.999	0.000
subject_racehispanic	0.941	0.224	4.202	0.000
subject_racewhite	0.826	0.224	3.686	0.000
subject_raceother	0.796	0.355	2.242	0.025
subject_raceunknown	-0.209	0.416	-0.502	0.616

Other Results			
Contraband found	Drugs	Weapons	Others
(Intercept)	-5.24*** (1.00)	-3.61*** (0.45)	-2.38*** (0.26)
Black	1.10 (1.01)	0.32 (0.46)	0.99*** (0.26)
Hispanic	1.21 (1.01)	0.36 (0.46)	1.04*** (0.26)
White	0.73 (1.01)	0.98* (0.46)	0.72** (0.26)
Other	0.97 (1.42)	0.47 (0.74)	0.87* (0.40)
Unknown	-11.33 (280.85)	0.04 (0.85)	-0.23 (0.53)

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

From table 7, black and Hispanic people are more likely to be frisked than white people. The estimated odd of being frisked for the black is 2.22 times the estimated odd for the white. This odd ratio for Hispanic people is 1.8. Asian people is the least likely to be frisked. From table 8 and table 9, contraband items is more likely to be found for Hispanic and black people. White people is more likely to be found with weapons and black and Hispanic people are more likely to be found with contraband items that are neither drugs or weapons.

Bayesian Modeling

Investigating the Hit Rate

The “hit rate,” defined here as the proportion of times an officer finds contraband given that a frisk has been performed, is a widely-used measure for assessing potentially-discriminatory policing. The hit rate can be thought of as a proxy for “evidence” when an officer decides whether to conduct a search or a frisk; a lower hit rate for a particular segment of the population can signal that an officer has a lower threshold of evidence when policing that population segment. In the following analysis, we examine the hit rate at the officer level. Because the analysis requires that officers have stopped all races under consideration, we restrict the analysis to only White, Black, and Hispanic subject races and to officers with 18 or more stops, corresponding to roughly the 90th percentile.

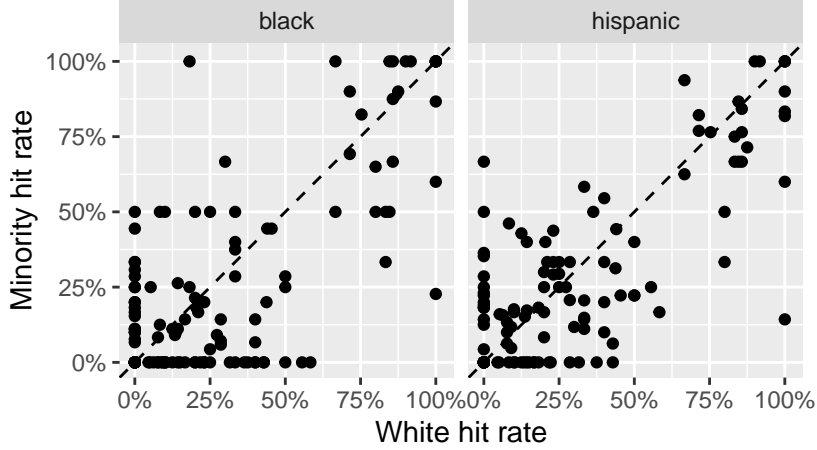


Figure 3: Hit rates for individual officers.

The above plot shows the hit rates for individual officers. An officer with an identical hit rate for white and minority subpopulations would be on the 45-degree line. Visually, it is difficult to determine a systematic trend, although it is clear that particular officers have hit rates that differ substantially by subpopulation. It should be noted that the hit rate is highly variable with small sample sizes.

We proceed using a Bayesian hierarchical model. Under this model, we treat individual officers as belonging to a population of players and we seek to model both the hit rates of the officers and the variation of this population. This permits *partial pooling*, by which individual hit rates are biased towards the population average by an amount determined by the estimate of the population. For each officer, we consider three hit rates, one each for white, Black, and Hispanic subpopulations. We accomplish this by fitting separate logistic mixed effects models for each race, each with a weakly informative Normal prior on the log-odds with mean -1.2 and standard deviation 1.

Specifically, let θ_{jr} be the hit rate for officer j and race r , y_{jr} be the number of hits, and K_{jr} the number of frisks. In the following, because we fit separate models, we assume for example $r = \text{black}$ and drop the r subscript. Assuming each officer's searches are independent Bernoulli trials

$$p(y_j|\theta_j) = \text{Binomial}(y_j|K_j, \theta_j)$$

We reparametrize the model in terms of the log-odds, α :

$$\alpha_j = \text{logit}(\theta_j) = \log \frac{\theta_j}{1 - \theta_j}$$

We set a weakly informative prior centered at $\alpha_j = -1.3$, corresponding to $\theta_j \approx 0.2$. The model is therefore

$$p(y_j|K_j, \alpha) = \text{Binomial}(y_j|K_j, \text{logit}^{-1}(\alpha_j))$$

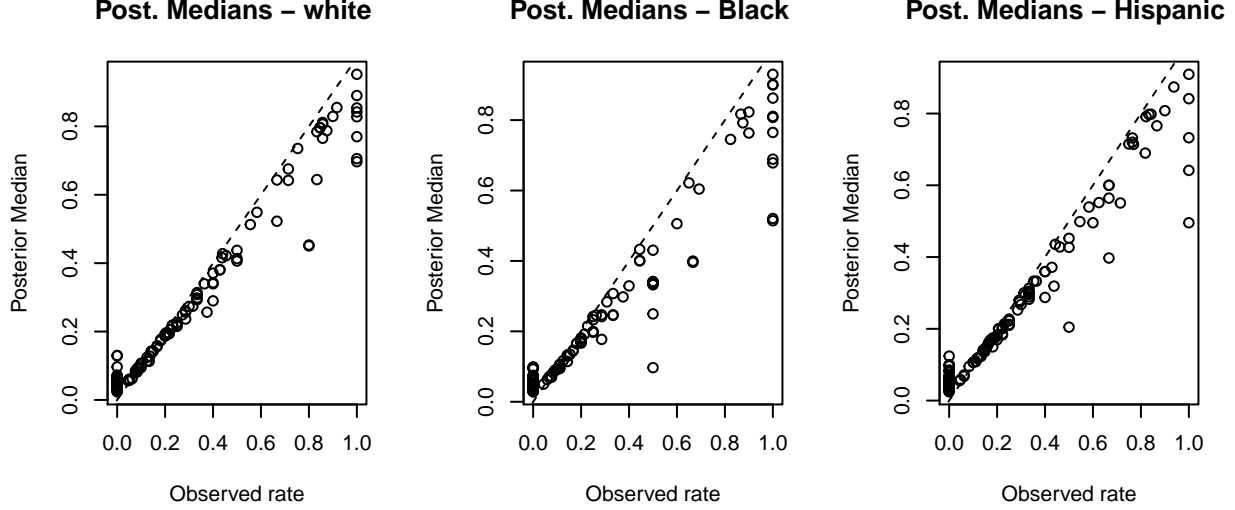
We proceed using `stan_glmr` and the default prior on the covariance matrix. The result includes a posterior for each officer; we may transform from the log-odds back to hit rate to obtain a posterior for the hit rate for each officer. We model each race separately, and so obtain three posteriors for each officer. Table 9 shows the posteriors for the first several officers (rows) for each of the three races (columns)

The following, the effects of partial pooling are evident: the posterior medians are biased towards the population average. Practically, this means that observed hit rates equal to zero have posterior medians that are small but positive, and perfect (or near-perfect) observed hit rates have somewhat smaller posterior

Table 9: Posterior intervals for three races. From left to right: white, Black, Hispanic

2.5%	50%	97.5%	2.5%	50%	97.5%	2.5%	50%	97.5%
0.003	0.055	0.317	0.021	0.198	0.617	0.065	0.214	0.450
0.021	0.095	0.248	0.002	0.033	0.203	0.047	0.162	0.359
0.003	0.056	0.324	0.111	0.242	0.420	0.037	0.174	0.423
0.128	0.523	0.900	0.057	0.521	0.959	0.665	0.874	0.974
0.054	0.188	0.415	0.032	0.332	0.833	0.082	0.268	0.556
0.114	0.380	0.722	0.002	0.040	0.261	0.010	0.071	0.237

medians.



Because part of this project is to “operationalize” fairness, we devised a measure by which the above posteriors can be converted into a rough “fairness score.” Because an officer that uses the same evidence threshold when deciding whether to frisk a subject regardless of race should have roughly equal hit rates for all three subpopulations, we reason that such an officer should have posterior medians that are close to each other for the three subpopulations. So, one can calculate a simple sum of squares statistic for each officer. Specifically, letting m_{jr} be the posterior median for officer j and race r , the sum of squares statistic S_j is

$$S_j = \sum_r (m_{jr} - \bar{m}_j)^2$$

where \bar{m}_j is the average of the three medians. Of course, this measure disregards all other information that could be gleaned from the posterior; an alternative might calculate the overlap between the posterior densities. However, we think this measure is relatively easy to understand and implement.

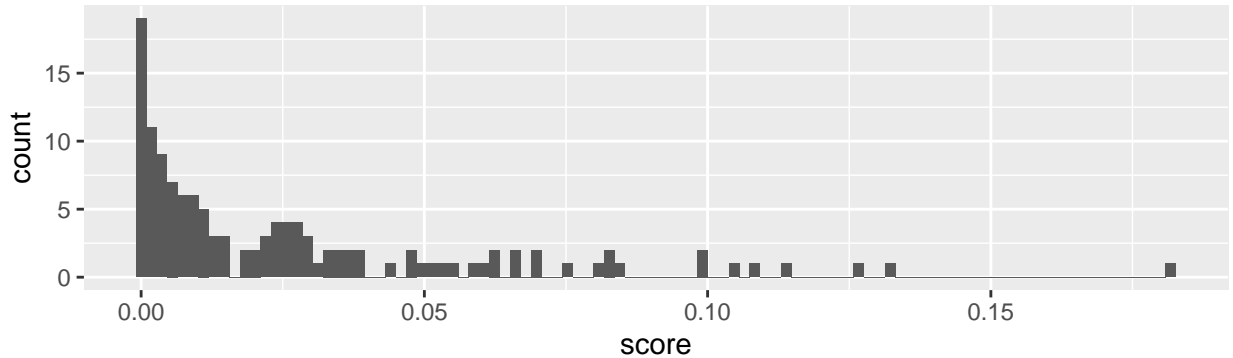


Figure 4: Fairness scores for the officers under consideration. Lower scores indicate hit rates are more similar.

Table 10: Highest 5 scores.

ID	obs.w	obs.b	obs.h	post.w	post.b	post.h	count.w	count.b	count.h
1504c3bc16	1.000	0.227	0.143	0.697	0.215	0.142	2	22	7
3392a495a3	0.400	0.000	0.545	0.371	0.046	0.499	10	5	11
50f70c6ecb	0.583	0.000	0.167	0.549	0.064	0.162	12	3	18
bab7c2acaf	0.833	0.333	0.750	0.644	0.247	0.715	7	3	20
dd9c1003d5	0.556	0.000	0.250	0.513	0.043	0.210	9	6	4

Conclusion

In this study, we evaluate the fairness of traffic stops during the past two decades through three tests, namely benchmark test, outcome test and veil of darkness test. We found that the racial disparity in policing exists and is present in different scales spatially and temporarily. We also explore the causal confounding issues through logistic regression, finding that black and Hispanic people are more likely to be frisked and found with contraband items that are neither drugs or weapons.

Through the investigation of the hit rate via Bayesian hierarchical modeling, we obtained posteriors for the hit rate for each officer in a subset of the data. Using the medians for these posteriors, we devised a “fairness score,” a tool we believe could be used to identify officers with racially disparate patterns of traffic stops.



Figure 5: FunkyStats: David Skril, Qiuyi Wu, Cuong Pham (from left to right)

References

Grogger, Jeffrey, and Greg Ridgeway. 2006. *Testing for Racial Profiling in Traffic Stops from Behind a Veil of Darkness*. Santa Monica, CA: American Statistical Association.

“Hierarchical Partial Pooling for Repeated Binary Trials.” n.d. <https://mc-stan.org/users/documentation/case-studies/pool-binary-trials.html>.

Pierson, Emma, Camelia Simoiu, Jan Overgoor, Sam Corbett-Davies, Daniel Jenson, Amy Shoemaker, Vignesh Ramachandran, et al. 2020. “A Large-Scale Analysis of Racial Disparities in Police Stops Across the United States.” *Nature Human Behaviour* 4 (7): 736–45. <https://doi.org/10.1038/s41562-020-0858-1>.

Simoiu, Camelia, Sam Corbett-Davies, and Sharad Goel. 2017. “THE Problem of Infra-Marginality in Outcome Tests for Discrimination.” *The Annals of Applied Statistics* 11 (3): 1193–1216. <http://www.jstor.org/stable/26362224>.